

リング結合型並列計算機のシステムアーキテクチャ

中條 拓伯* 和田 耕一** 金田 悠紀夫*

* 神戸大学 工学部 システム工学科

** 筑波大学 電子情報工学系

リング結合型並列計算機は疎結合でありながら、ネットワーク仮想共有メモリなど、密結合マルチプロセッサのもつ利点をも有効に活かすためのマシンであり、1024台までの拡張を目指して現在開発中である。ここではシステム全体のアーキテクチャと処理結果の回収のためのホストコンピュータと各要素プロセッサ(PE)との通信形態と割込みについて述べ、ネットワークの性能を低下させないためのフロー制御について説明する。さらに、B-treeをPE上に構築し、並列に検索、挿入、削除する方法について解説する。

S y s t e m A r c h i t e c t u r e o f

R i n g - C o n n e c t e d P a r a l l e l C o m p u t e r

H i r o n o r i N a k a j o * K o i c h i W a d a ** Y u k i o K a n e d a *

* D e p a r t m e n t o f S y s t e m s E n g i n e e r i n g ,
 F a c u l t y o f E n g i n e e r i n g , K o b e U n i v e r s i t y
 R o k k o d a i - c h o , n a d a - k u , K o b e , 6 5 7 J a p a n ,

** I n s t i t u t e o f I n f o r m a t i o n S c i e n c e a n d
 E l e c t r o n i c s , U n i v e r s i t y o f T s u k u b a
 1 - 1 - 1 , T e n n o d a i , S a k u r a m u r a , N i h a r i g u n ,
 I b a r a k i k e n , 3 0 5 J a p a n

Though our Ring-Connected Parallel Computer is a loosely-coupled multi-processor, equipped with network virtual shared memory, it has characteristics of a tightly-coupled multi-processor. This parallel machine is now being developed with 32 processors, but, in the future, we aim the expansion to 1024 processors.

We describe the system overall configuration, communication method between Host computer and each processing element(PE)s to collect answers processed in each PE, interrupt mechanism, and flow control not to deteriorate network performance.

Moreover, we explain the parallel search, insertion, and deletion using B-tree constructed in each PE.

1. はじめに

リング結合型並列計算機はデータベース、知識ベースの並列処理、さらにネットワーク上での仮想共有メモリなどに関する研究を行うための実験システムとして現在、開発が進んでいる。

ここでは、本システムの構成と割り込み機構について述べ、さらにネットワーク内の輻輳を回避するためのフロー制御について説明する。

そして、本システムにおいてB-treeによる検索、挿入、削除を並列に行う方式を提案し、その可能性について考える。

2. システムアーキテクチャ

図1にリング結合型並列計算機のシステム全体の構成を示す。隣接する要素プロセッサ (PE) はストリーム・コントローラ (SC) とよぶ通信専用ハードウェアによってリング状に結合されている。PE間はSCにより、高速にパケット通信が行われ、負荷分散やいろいろなメッセージ通信に用いられる。さらに、ホストコンピュータから各PEはブロードキャスト・バスで接続されており、ホストからPEへ同一のデータが同時に転送 (放送) される。このとき、すべてのPEに対して、またはいくつかのPEに対して (受信者選択型マルチキャスト)、さらに個々のPEに対してそれぞれ選択して転送できる^{[1][2]}。

表1に各PEのハードウェア仕様を示す。

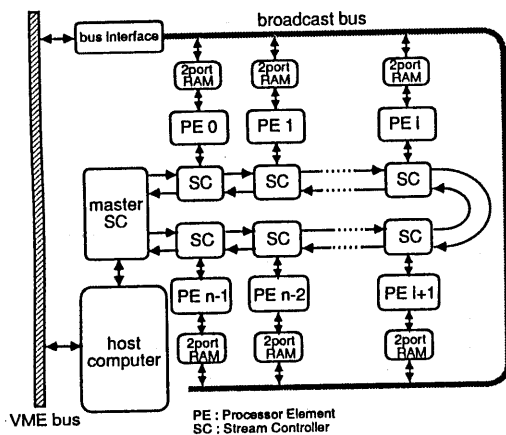


図1 リング結合型並列計算機のシステム構成

規模:	最大 1024プロセッサ
接続形態:	リング+ブロードキャスト
CPU:	MC68020 (16MHz)
ローカルメモリ:	ROM 128Kbyte
	DRAM 2Mbyte
放送メモリ:	SRAM 512Kbyte

表1 各PEのハードウェア仕様

2.1 システムのメモリマップ

図2にホストコンピュータと各PEのメモリマップを示す。各PEは放送領域として512KバイトのSRAMをインストールしているが、この領域は、図2に示すようにホストのメモリマップ上にもマッピングされており、ホストは1024台のPEから構成した場合、最大512Mバイトの実メモリをインストールすることになる。どのPEに対して転送するか、すなわち、どのようにマルチキャストを行うかは、4Gという広大なメモリ空間を有するMC68020のアドレスバスを有効に用いて決定される^[1]。

各PEは2MバイトのDRAMをインストールしているが、その一部をネットワーク仮想メモリとして用い、他のPEと共有され、仮想的に個々のPEは最大512Mバイトのネットワーク仮想共有メモリを有することになる。これについては後述する。

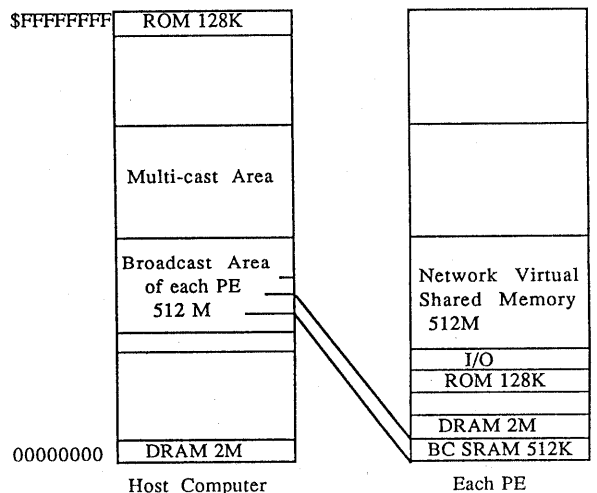


図2 メモリマップ

2. 2 割込み

マルチプロセッサシステムにおいて、ホストコンピュータが各PEへの処理開始の合図や、各PEから処理の終了の報告を行うときに割込みが有効に用いられる。ここでは、ホストコンピュータから各PEへ、また各PEからホストへの2種類の割込みについて説明し、その有効性について述べる。

(1) ホストコンピュータから各PEへの割込み

データベース、知識ベースなどの大量のデータを分割し、各PEに分配する場合を考える。PE側でローディングの終了を知る方法としてEnd of Data (EOD)を表す特殊なデータを用いる。ここで、PEが他に行うべき処理がなければ、EODをポーリングすることによって処理の開始時を知ることができるが、他の処理を行っている場合は割込みが有効に用いられる。

具体的には、前に述べたマルチキャストを用いて、ホストから個々のPE、または複数のPEに対して割込みをかけることができる。つまり、ホストが各PEの放送領域の最終アドレスにデータを書き込んだ場合に割込み信号を発するが、そのときのアドレスによってどのPEに対して割込みをかけるかを容易に決定できる。

割込みをかけられたPEは、そのときに優先順位より高い割込みがかかっていなければ、ホストからの割込みに対するアクナリッジ信号をオープンコレクタバスにH (High)として出力する。PEがたとえ一つでもアクナリッジ信号をL (Low)で返すと割込みは受け付けられず、アクナリッジがHになるまで待たされることになる。

(2) 各PEからホストコンピュータへの割込み

それぞれのPEにおいて処理した結果の回収について考える。各PEの放送領域はホストコンピュータのメモリマップ上にマッピングされているので、PEはその領域に処理結果を書き込み、ホストが読みに行くことで容易に結果を回収できる。

例えば、データベースをPEに分配し、検索データを放送して、検索を行う場合を考える。各PEの放送領域にステータスフラグを用意し、PEが処理を終了したかどうか、解を持つかどうかなどの情報を蓄える。このとき、

a. 多数のPEが解を持つ場合

ホストは各PEのステータスフラグを順に読み込み、処理が終了していれば、結果を読み込んで回収する。PEがまだ実行中であれば、そのPEのIDをキュー (process queue) に書き込み、次のPEのステータスフラグを読み込む。以上の操作をすべてのPEに対して行い、1通りの処理が終了すれば、キューに格納されたPEに対して、再びステータスフラグを読み込み、キューが空になるまで続ける。以上の操作を図3に示す。

```
for (i = 0; i < n; i++) {
    if (status_flag(i) != finish)
        process_enqueue(i);
    else
        collect_result(i);
}

do {
    id = process_dequeue();
    if (status_flag(id) != finish)
        process_enqueue(id);
    else
        collect_result(id);
} while (process_queue == empty);
```

図3 複数PEからの結果回収方法

b. 少数のPEが解を持つ場合

この場合、ホストが各PEのステータスフラグを順次チェックするのは効率が悪く、それよりPE側から解を持っていることをホストに知らせる方が有効と考えられる。そこで、PEからホストへの割込みが用いられる。

VMEバスなど、単一バスに複数のプロセッサが接続されているようなマルチプロセッサシステムでは、個々のボードからの割り込みはディージ・チェーンによって排他制御され、どのボードからの割り込みかを判断するには、各ボードから送出されるベクタを用いる。

しかし、MC68020の場合、ベクタ・テーブルは256個に限られており、さらにMPU自身が64個のベクタを、アドレス・エラーやゼロ除算などのために定義している。したがって、ベクタによってインタラプタを判断するには、192台までが限度であり、1024台までの拡張を考えている本システムでは使用は困難である。さらに、ディージ・チェーンを用いた場合、その伝搬時間も見過ごすことのできない問題となる。

そこで、本システムでは割り込みのための割り込みIDバスとよばれる特別なバスを設け、他のPEからの割り込みの競合に対して排他制御を行っている。ホストは1024台のPEの中で、どのPEからの割り込み要求かを容易に判断でき、それぞれの要求に対して応答する。PEからの割り込み発生のためのハードウェア・ブロック図を図4に示す。

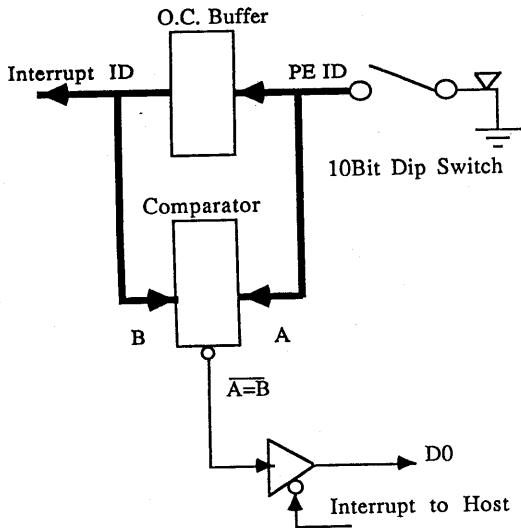


図4 PEからの割り込みのブロック図

PEからホストへの割り込みシーケンスについて説明する。

① PEが割り込みをかける場合、あるポートを読みに行く。

② そのポートのアドレスをデコードして、自身のIDをインタラプト・ID (INT・ID) バスに乗せる。

③ INT・IDバスはオープンコレクタ出力で他のPEのINT・IDバスと接続されており、他に割り込みをかけているPEがなければ、自身のIDとINT・IDとが一致し、その結果をデータバスの最下位ビット (D0) に乗せる。①でポートを読みに行くと言うのは、実はこのD0を読みに行くことであり、もし0ならば他に競合することがなく、ホストに割り込みを受け入れられたと判断でき、1ならば他のPEと競合し、割り込みを受け入れられなかったと判断できる。①から③までは、MC68020の単なるリードの1サイクルであり、ノー・ウェイトもしくは1・ウェイトで動作する。

④ ホストは割り込みを受けつけた後、このINT・IDを読みに行くことで、どのPEからの割り込みを容易に判断できる。INT・IDは10ビット用意してあるので、ホストは1024台までのPEからの割り込みを判断することができる。

⑤ クリティカルなタイミングで、2つ以上のPEが同時に割り込みをかけた場合、ホストへの割り込みはかからない。

PEからホストへの割り込み操作を図5に示す。

```
do {
    check = interrupt_to_host();
    if (LSB(check) == 1)
        do_another();
} while (LSB(check) == 0);
```

図5 PEからホストへの割り込み

3. フロー制御とデッドロックの回避

複数のノードからなるネットワーク上の問題として輻輳 (congestion) の発生がある。これはメッセージやデータが一箇所、または数カ所に集中するために起こる。このため、メッセージの伝達速度は急速に低下し、最終的にはデッドロックに陥る恐れがある。

一般的なネットワーク上では次の2種類のデッドロックが起こり得る^[6]。

a. プロトコル・デッドロック

受信確認応答や送信許可など、ある種の制御信号のような論理的な資源が原因となって発生する。

b. バッファ・デッドロック

メッセージ・バッファなどの物理的な資源が原因となって発生する。

本システムでは、ネットワークの性格上、バッファ・デッドロックが起こる可能性が高いと考えられる。ネットワークにおいて輻輳を防ぐためのフロー制御には、ネットワーク内でただ1つのノードが制御する集中型フロー制御と、各ノードで制御する分散型フロー制御がある。本システムではブロードキャスト・バスを用いて集中的に管理することも可能であるが、ホストの負荷を考慮して、分散型を選択する。一般的なフロー制御には次の3つの方法がある^[6]。

- ・適応ルーティング
- ・エンドツーエンド・フロー制御
- ・イザリズミック・フロー制御

ルーティングを不要にするためにリング構造を採用した本システムでは、適応ルーティングは無関係である。エンドツーエンド・フロー制御は受信先の状態にもとづいて送信側がメッセージ・トラフィックを発生させる割合を規制するものであるが、本システムでは、受信側の状態を知るようなメッセージのやりとりを行うようなことは考えていない。イザリズミック・フロー制御はネットワーク内のメッセージ総数を一定値に制限することで輻輳を制御するものである^[6]。

本システムでは、このイザリズミック・フロー制御を採用する。具体的には、ネットワーク内に一定数のパーミット (入場券) が発行され、このパーミットを得たパケットだけがネットワーク内で転送される。このため、ネットワーク内のパケット数はこのパーミット数以下に保たれ、地域的に多少の輻輳が生じるが、ネットワーク全体に大きな影響を及ぼすことはない。パーミットを得てネットワーク内に入ったパケット、またはすでに伝送されているパケットを内部パケットとよび、それ以外を外部パケットとよぶ。外部パケットが伝送されるためには、そのパケットを持つノードは他のノードからパーミットが配送されるのを待たなければならない。内部パケットには次の2種類のパケットが存在する。

(1) データパケット

通常のメッセージあるいは制御情報を含むパケットで、通常1個のパーミットを有し、転送先のノードでそのパーミットを提供する。本システムでは、データパケットをさらにエンドツーエンド (end-to-end) ・データパケットとブロードキャスト (broadcast) ・データパケットの2つに分類する。前者はPE間における1対1通信に用いられ、後者は1対多通信に用いられる。このブロードキャスト・データパケットを用いることで、すべてのPEは他のPEへ放送を行うことができる。

(2) 空パケット

このパケットはパーミットの配送のみに用いられ、ある最大数以下のパーミットを一時に運ぶ。ノードにおいてパーミットがあふれた場合、余分のパーミットは、この空パケットを用いることで他のノードに配送される。空パケットはメッセージを含んでいないためネットワーク上のオーバーヘッドとなる。

このようにネットワーク内のパーミット数を一定値以下にすることで、ネットワーク全体のパケット数を制限できる。本システムでの実際の制御は以下の手順で行う。

a. 空パケットを受信する場合

① パーミット・バッファがフルならば、そのパーミットをエンドツーエンド空パケットとして、ネットワークの隣接PEに転送する。(図6-1)

② パーミット・バッファがフルでないならば、その空パケットを取り込み、パーミット・バッファにパーミットを追加する。(図6-2)

b. エンドツーエンド・データパケットを受信する場合

① パーミット・バッファがフルならば、データパケットを取り込んだ後、空パケットをネットワーク上に発信する。(図7-1)

② パーミット・バッファがフルでないならば、データパケットを取り込んだ後、パーミット・バッファにパーミットを1つ追加する。(図7-2)

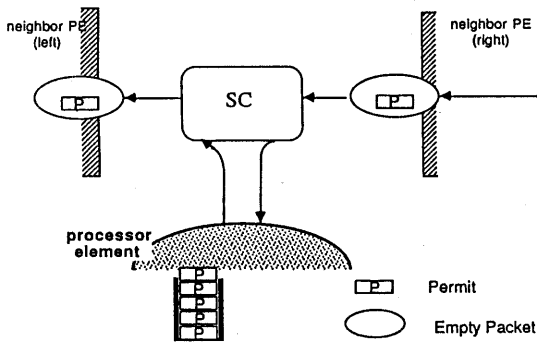


図6-1 空パケットの受信
(パーミット・バッファがフル)

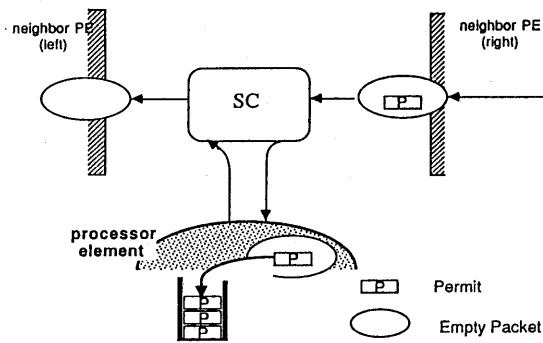


図6-2 空パケットの受信
(パーミット・バッファがフルでない)

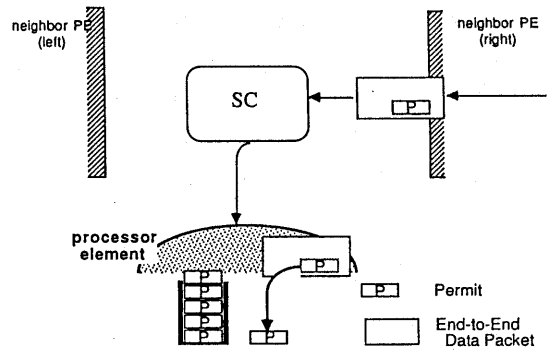


図7-1 エンドツーエンド・データパケット
の受信 (パーミット・バッファがフル)

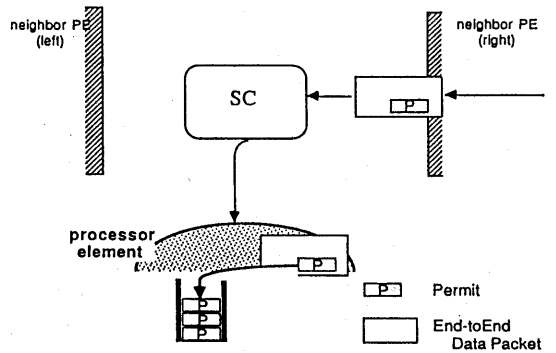


図7-2 エンドツーエンド・データパケット
の受信 (パーミット・バッファがフルでない)

c. ブロードキャスト・データパケットを受信する場合 (図8)

ブロードキャスト・データパケットはネットワークを巡回し、それぞれのPEに取り込まれ、一回りした後に発信ノードによって回収される。このときブロードキャスト・データパケットはエンドツーエンド・データパケットと同様にパーミットを含むが、パーミットの配送をそれぞれのPEに対して行えば、パーミット数が著しく増加する。そこで、ブロードキャスト・データパケットはパーミットの配送は行わないことにする。すなわち、各PEはブロードキャスト・データを受け取るが、パーミットは受け取らない。

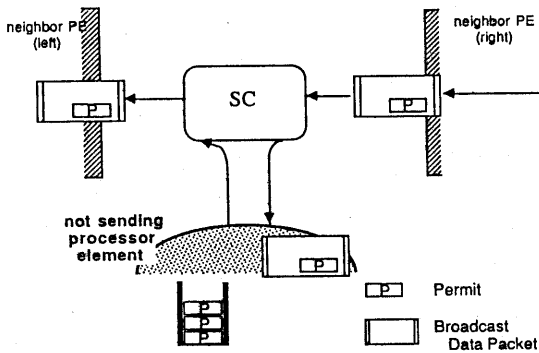


図8 ブロードキャスト・データパケットの受信

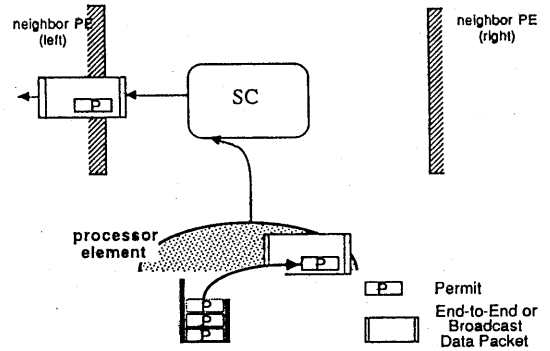


図9-1 エンドツーエンド・データパケットの送信 (パーミット・バッファが空でない)

d. エンドツーエンド・データパケットの送信

(1) パーミット・バッファが空でないならば、エンドツーエンド・データパケットを送信した後、パーミット・バッファからパーミットを1つ削除する。(図9-1)

(2) パーミット・バッファが空ならば、空パケットが隣接PEから転送されパーミットが配送されてくるまで送信を待つ。(図9-2)

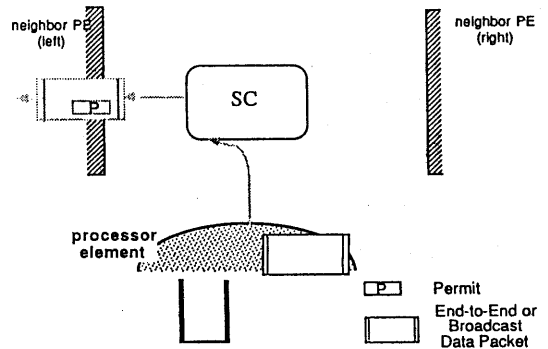


図9-2 エンドツーエンド・データパケットの送信 (パーミット・バッファが空)

e. ブロードキャスト・データパケットの送信

(1) パーミット・バッファが空でないなら、ブロードキャスト・データパケットを送信した後、パーミット・バッファからパーミットを1つ削除する。

(2) パーミット・バッファが空ならば、空パケットが転送されてくるまで送信を待つ。

ネットワーク上の総パーミット数についてであるが、パーミットが永遠に得られないという飢餓状態を回避するためにネットワーク上には、絶えずいくつかの空パケットを保持する必要がある。パーミット数を多く取れば、パーミット・バッファが空となるPEは少なくなるが、空パケットの送受信のためのオーバーヘッドが大きくなる。パーミット数を少なくすると、パーミットを持たないPEが増加し、データパケットの転送効率は低下する。したがって、転送スピードを考慮して最適なパーミットの分配を行わなければならない。

4. ネットワーク仮想共有メモリ

本システムでは、高速なパケット通信を用いて、各PEのメモリの一部を他のPEと共有することで、仮想的に大容量のメモリを有するというネットワーク仮想共有メモリの実現を構想している。共有メモリの有効性については、比較的少数のPEからなる商用機の成功が示す通りであるが、PE数の増加に対しては悲観的である。そこで、ネットワークを用いたシステムによる仮想共有メモリが今後有利になると考えられる。図2のメモリマップに示したように、各PEは512Mバイトの仮想共有メモリを持つ。現在は、データの無矛盾化プロトコルについて詳細を検討中である。

5. B-treeの並列実行方式について

データベースに蓄積されるレコードの数が増加するに従って、その中から特定のデータを取り出すことが困難になってくる。あるキーを与えて、それと一致するキーの値を持つレコードをファイルの中から検索する方法の代表的なものに、ハッシングとインデックスがある。データを検索する場合、ハッシュを用いればO(1)で検索できるが、検索するキーをある範囲で指定した場合にはインデックスの方が有効である。また、データの追加に関して、ハッシングでは追加データが増加するとハッシュテーブルを再編成する必要が生じ、その処理とテーブルに消費するメモリ容量が問題になってくる。ここでは、インデックスの中でも特にB-treeについて、本システムでの並列実行の可能性について述べる。

まず、ホストコンピュータの2次記憶装置にあるデータベースを各PEへ分配する。この場合に次に示す2通りの方法がある。分配を行った後、各PE内でB-treeの構築を行う。

(1) キーの値において、どのPEがどの範囲のデータを処理するかを決定しておき、ブロードキャスト・バスを用いてデータを分配する。

(2) ホストはデータベースを等量に分割し、それぞれを各PEに分配する。

(1)の分配はブロードキャスト・バスを用いれば比較的容易に実現できる。すなわち、個々のPEの放送領域がホストのメモリマップ上にマッピングされているので、キーの値から送るべき領域の計算(アドレス変換)を行うのに算術計算だけでよく、複雑な条件処理を行う必要はない。ホストがあるデータを検索、挿入、削除を行うときに、そのデータのキーの値から、どのPE内のtreeにあるかを判断できるので、多数のデータを検索する場合、データの検索を行うPEに随時転送し、並列に処理を進めることが可能である。ただし、キーの値によって、PE間でデータ数のばらつきが生じる。この場合、隣接PE間でデータ転送を行い、データ数になるべく等しくなるようにする必要がある。

(2)の場合は、(1)よりもデータの分配が高速に行えるが、データの検索は、検索データを放送することでできる。データの挿入を行う場合も、すべてのPEにそのデータを放送し、どのtreeにも含まれないことがわかれば、データ数の少ないPE内のtreeに挿入すればよい。各PEの持つデータ数は、ほぼ等量に保てるので負荷分散の点では(1)より優れていると思われる。

6. おわりに

現在は32台からなるプロトタイプを製作中である。B-treeの並列処理に関する問題は既存の密結合マシンにおいて評価をおこなっている。今後は、ネットワーク仮想共有メモリについても、その詳細を明確にしていく予定である。

【参考文献】

- [1]中條, 和田, 金田, 他: "論理型並列データベースマシンに関する研究", 信学技報, CPSY88-29, (1988)
- [2]中條, 和田, 金田, 他: "並列知識ベースマシンのアーキテクチャ", 第37回情報処理学会全国大会, 3N-2, (1988)
- [3]中條, 和田, 金田: "並列知識ベースマシンのハードウェア構成", 並列処理シンポジウムJSP'89, PP.147-154(1989)
- [4]富田: 並列計算機構成論, 昭晃堂(1986)
- [5]Hsiao, D.K.: "The Impact of the Interconnecting Network on Parallel Database Computers", Database Machines and Knowledge Base Machines, Kluwer Academic Publishers(1988)
- [6]Ahuja, V.A.: "Design and Analysis of Computer Communication Networks", McGRAW-HILL(1982)